```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```python
house=pd.read_csv("C:/Users/TERIAK-JB/Desktop/GOMYCODE/kc_house_data.csv",encoding="iso-8859-1")
```

```python
house.head(10)
```

| | id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | ... | grade | sqft_abov |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7129300520 | 20141013T000000 | 221900.0 | 3 | 1.00 | 1180 | 5650 | 1.0 | 0 | 0 | ... | 7 | 118 |
| 1 | 6414100192 | 20141209T000000 | 538000.0 | 3 | 2.25 | 2570 | 7242 | 2.0 | 0 | 0 | ... | 7 | 217 |
| 2 | 5631500400 | 20150225T000000 | 180000.0 | 2 | 1.00 | 770 | 10000 | 1.0 | 0 | 0 | ... | 6 | 77 |
| 3 | 2487200875 | 20141209T000000 | 604000.0 | 4 | 3.00 | 1960 | 5000 | 1.0 | 0 | 0 | ... | 7 | 105 |
| 4 | 1954400510 | 20150218T000000 | 510000.0 | 3 | 2.00 | 1680 | 8080 | 1.0 | 0 | 0 | ... | 8 | 168 |
| 5 | 7237550310 | 20140512T000000 | 1225000.0 | 4 | 4.50 | 5420 | 101930 | 1.0 | 0 | 0 | ... | 11 | 389 |
| 6 | 1321400060 | 20140627T000000 | 257500.0 | 3 | 2.25 | 1715 | 6819 | 2.0 | 0 | 0 | ... | 7 | 171 |
| 7 | 2008000270 | 20150115T000000 | 291850.0 | 3 | 1.50 | 1060 | 9711 | 1.0 | 0 | 0 | ... | 7 | 106 |
| 8 | 2414600126 | 20150415T000000 | 229500.0 | 3 | 1.00 | 1780 | 7470 | 1.0 | 0 | 0 | ... | 7 | 105 |
| 9 | 3793500160 | 20150312T000000 | 323000.0 | 3 | 2.50 | 1890 | 6560 | 2.0 | 0 | 0 | ... | 7 | 189 |

10 rows × 21 columns

◄ |                                                                    | ►

```python
house.shape
```

```
(21613, 21)
```

```python
house.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21613 entries, 0 to 21612
Data columns (total 21 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   id           21613 non-null  int64
 1   date         21613 non-null  object
 2   price        21613 non-null  float64
 3   bedrooms     21613 non-null  int64
 4   bathrooms    21613 non-null  float64
 5   sqft_living  21613 non-null  int64
```

```
 6   sqft_lot       21613 non-null  int64
 7   floors         21613 non-null  float64
 8   waterfront     21613 non-null  int64
 9   view           21613 non-null  int64
 10  condition      21613 non-null  int64
 11  grade          21613 non-null  int64
 12  sqft_above     21613 non-null  int64
 13  sqft_basement  21613 non-null  int64
 14  yr_built       21613 non-null  int64
 15  yr_renovated   21613 non-null  int64
 16  zipcode        21613 non-null  int64
 17  lat            21613 non-null  float64
 18  long           21613 non-null  float64
 19  sqft_living15  21613 non-null  int64
 20  sqft_lot15     21613 non-null  int64
dtypes: float64(5), int64(15), object(1)
memory usage: 3.5+ MB
```

In [6]:

```
house.drop("yr_renovated",1,inplace=True)
```

In [7]:

```
house["bathrooms"].apply(np.round)
house["bathrooms"]=house["bathrooms"].apply(np.int)
```

In [8]:

```
house=house.drop(house[(house.bedrooms>10)].index)
```

In [9]:

```
house
```

Out[9]:

| | id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | condition | grade |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7129300520 | 20141013T000000 | 221900.0 | 3 | 1 | 1180 | 5650 | 1.0 | 0 | 0 | 3 | 7 |
| 1 | 6414100192 | 20141209T000000 | 538000.0 | 3 | 2 | 2570 | 7242 | 2.0 | 0 | 0 | 3 | 7 |
| 2 | 5631500400 | 20150225T000000 | 180000.0 | 2 | 1 | 770 | 10000 | 1.0 | 0 | 0 | 3 | 6 |
| 3 | 2487200875 | 20141209T000000 | 604000.0 | 4 | 3 | 1960 | 5000 | 1.0 | 0 | 0 | 5 | 7 |
| 4 | 1954400510 | 20150218T000000 | 510000.0 | 3 | 2 | 1680 | 8080 | 1.0 | 0 | 0 | 3 | 8 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 21608 | 263000018 | 20140521T000000 | 360000.0 | 3 | 2 | 1530 | 1131 | 3.0 | 0 | 0 | 3 | 8 |
| 21609 | 6600060120 | 20150223T000000 | 400000.0 | 4 | 2 | 2310 | 5813 | 2.0 | 0 | 0 | 3 | 8 |
| 21610 | 1523300141 | 20140623T000000 | 402101.0 | 2 | 0 | 1020 | 1350 | 2.0 | 0 | 0 | 3 | 7 |
| 21611 | 291310100 | 20150116T000000 | 400000.0 | 3 | 2 | 1600 | 2388 | 2.0 | 0 | 0 | 3 | 8 |
| 21612 | 1523300157 | 20141015T000000 | 325000.0 | 2 | 0 | 1020 | 1076 | 2.0 | 0 | 0 | 3 | 7 |

21611 rows × 20 columns

In [10]:

```
house["floors"].apply(np.round)
house["floors"]=house["floors"].apply(np.int)
```

```
house
```

Out[64]:

| | id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | condition | grade | sqft_above | sqft_ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7129300520 | 20141013T000000 | 221900.0 | 3 | 1 | 1180 | 5650 | 1 | 3 | 7 | 1180 | |
| 1 | 6414100192 | 20141209T000000 | 538000.0 | 3 | 2 | 2570 | 7242 | 2 | 3 | 7 | 2170 | |
| 2 | 5631500400 | 20150225T000000 | 180000.0 | 2 | 1 | 770 | 10000 | 1 | 3 | 6 | 770 | |
| 3 | 2487200875 | 20141209T000000 | 604000.0 | 4 | 3 | 1960 | 5000 | 1 | 5 | 7 | 1050 | |
| 4 | 1954400510 | 20150218T000000 | 510000.0 | 3 | 2 | 1680 | 8080 | 1 | 3 | 8 | 1680 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 21608 | 263000018 | 20140521T000000 | 360000.0 | 3 | 2 | 1530 | 1131 | 3 | 3 | 8 | 1530 | |
| 21609 | 6600060120 | 20150223T000000 | 400000.0 | 4 | 2 | 2310 | 5813 | 2 | 3 | 8 | 2310 | |
| 21610 | 1523300141 | 20140623T000000 | 402101.0 | 2 | 0 | 1020 | 1350 | 2 | 3 | 7 | 1020 | |
| 21611 | 291310100 | 20150116T000000 | 400000.0 | 3 | 2 | 1600 | 2388 | 2 | 3 | 8 | 1600 | |
| 21612 | 1523300157 | 20141015T000000 | 325000.0 | 2 | 0 | 1020 | 1076 | 2 | 3 | 7 | 1020 | |

21268 rows × 18 columns

```
house["waterfront"].value_counts()
```

Out[12]:

```
0    21448
1      163
Name: waterfront, dtype: int64
```

```
house.drop("waterfront",1,inplace=True)
```

```
house["view"].value_counts()
```

Out[14]:

```
0    19487
2      963
3      510
1      332
4      319
Name: view, dtype: int64
```

```
house.drop("view",1,inplace=True)
```

```
def plot_correlation_map( df ):
```

```
    corr = df.corr()

    s , ax = plt.subplots( figsize =( 30 , 30 ) )

    cmap = sns.diverging_palette( 250 , 5 , as_cmap = True )

    s = sns.heatmap(

        corr,

        cmap = cmap,

        square=True,

        cbar_kws={ 'shrink' : .50 },

        ax=ax,

        annot = True,

        annot_kws = { 'fontsize' : 15 }

        )
```
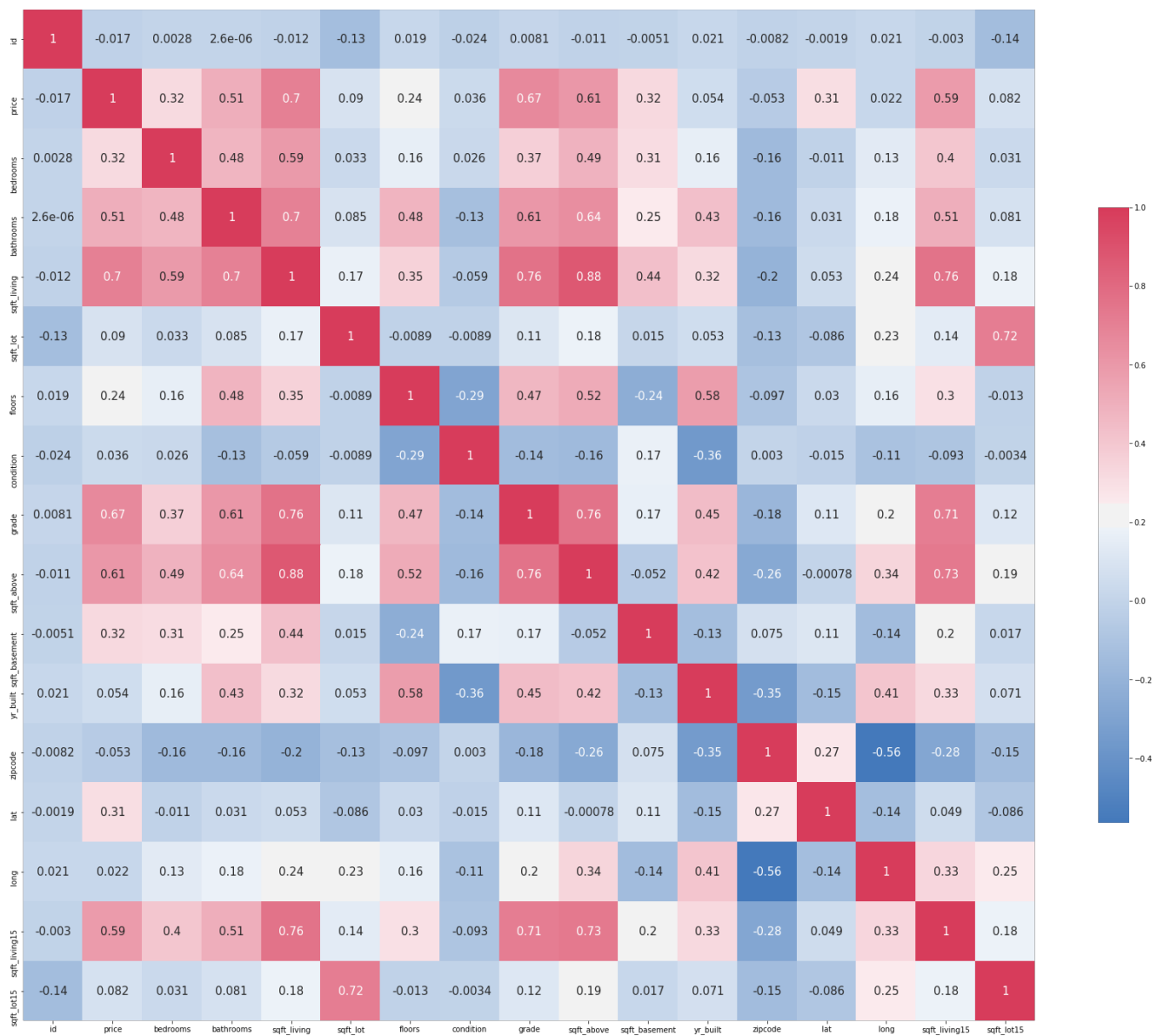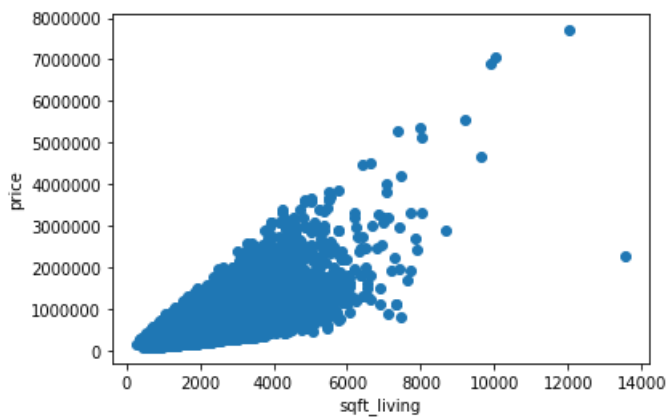
In [17]:

```
plot_correlation_map( house )
```

```python
plt.scatter(house["sqft_living"],house["price"])
plt.xlabel("sqft_living")
plt.ylabel("price")
```

Out[18]:

```
Text(0, 0.5, 'price')
```



In [19]:

```python
house["sqft_living"].value_counts()
```

Out[19]:

```
1300    138
1400    135
1440    133
1010    129
1660    129
        ...
3001      1
4970      1
2905      1
2793      1
1975      1
Name: sqft_living, Length: 1038, dtype: int64
```

In [20]:

```python
house["grade"].value_counts()
```

Out[20]:

```
7     8979
8     6068
9     2615
6     2038
10    1134
11     399
5      242
12      90
4       29
13      13
3        3
1        1
Name: grade, dtype: int64
```

In [21]:

```python
house=house.drop(house[(house.grade<6)].index)
```

In [22]:

```python
house.pivot_table('bathrooms', index='bedrooms', columns='grade', aggfunc="sum").plot()
```

```
house.pivot_table('bathrooms', index='bedrooms',columns='grade',aggfunc='sum').plot()
```

Out[22]:

<matplotlib.axes._subplots.AxesSubplot at 0xce3f508>



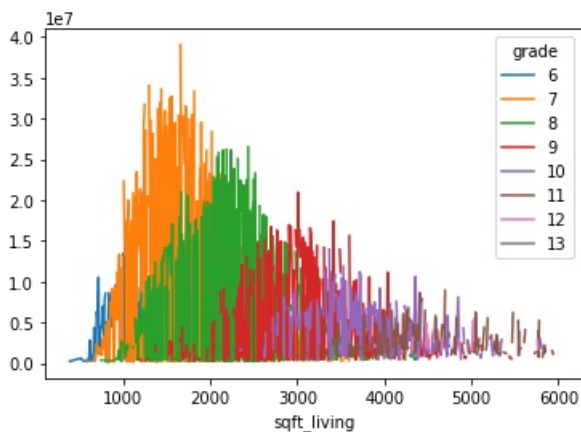In [81]:

```
house.pivot_table('price', index='sqft_living',columns='grade',aggfunc="sum").plot()
```

Out[81]:

<matplotlib.axes._subplots.AxesSubplot at 0xbb18108>



In [80]:

```
house=house.drop(house[(house.price>300000)&(house.sqft_living>6000)].index)
house.head()
```

Out[80]:

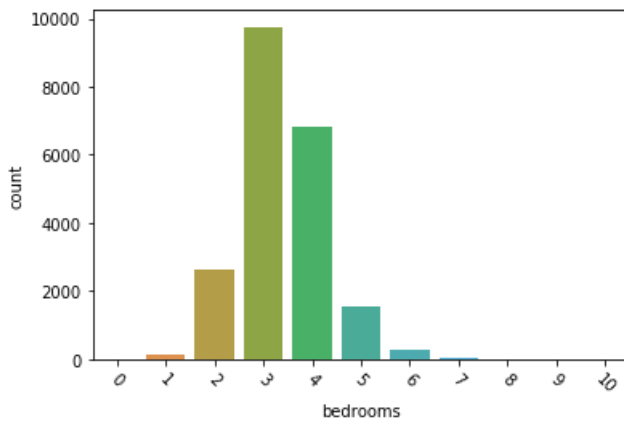| | id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | condition | grade | sqft_above | sqft_base |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7129300520 | 20141013T000000 | 221900.0 | 3 | 1 | 1180 | 5650 | 1 | 3 | 7 | 1180 | |
| 1 | 6414100192 | 20141209T000000 | 538000.0 | 3 | 2 | 2570 | 7242 | 2 | 3 | 7 | 2170 | |
| 2 | 5631500400 | 20150225T000000 | 180000.0 | 2 | 1 | 770 | 10000 | 1 | 3 | 6 | 770 | |
| 3 | 2487200875 | 20141209T000000 | 604000.0 | 4 | 3 | 1960 | 5000 | 1 | 5 | 7 | 1050 | |
| 4 | 1954400510 | 20150218T000000 | 510000.0 | 3 | 2 | 1680 | 8080 | 1 | 3 | 8 | 1680 | |

In [27]:

```
sns.countplot(x="bedrooms",data=house)
```

```
plt.xticks(rotation=-45)
```

Out[27]:

```
(array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10]),
 <a list of 11 Text xticklabel objects>)
```
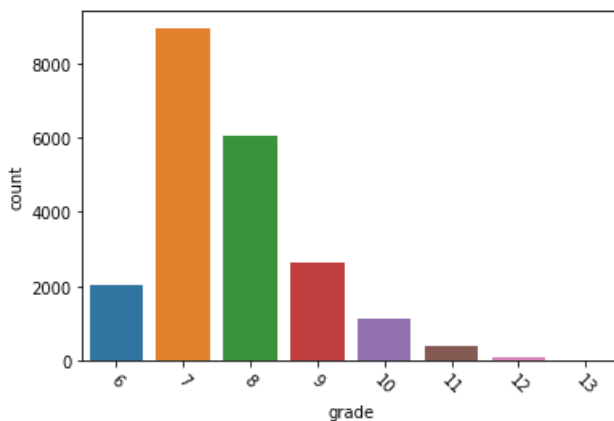


In [29]:

```
sns.countplot(x="grade",data=house)
plt.xticks(rotation=-45)
```

Out[29]:

```
(array([0, 1, 2, 3, 4, 5, 6, 7]), <a list of 8 Text xticklabel objects>)
```



Regression linéaire

In [30]:

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn import metrics
```

In [31]:

```python
x=house["price"]
y=house["sqft_living"]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20,random_state=30)
```

In [48]:

```
#extract x and y from our data
```

```
#extract x and y from our data
x=house["sqft_living"].values[:,np.newaxis]
y=house["price"].values


x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.35,random_state=40) #splitting data
with test size of 35%


model=LinearRegression()    #build linear regression model
model.fit(x_train,y_train)  #fitting the training data
predicted=model.predict(x_test) #testing our model's performance
print("MSE", mean_squared_error(y_test,predicted))
print("R squared", metrics.r2_score(y_test,predicted))
```
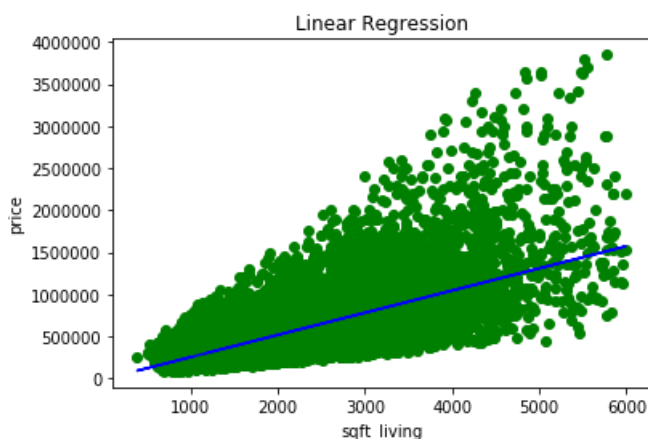
```
MSE 60164508415.7486
R squared 0.4680860370228618
```

In [33]:

```
plt.scatter(x,y,color="g")
plt.title("Linear Regression")
plt.ylabel("price")
plt.xlabel("sqft_living")
plt.plot(x,model.predict(x),color="b")
plt.show()
```



Regression Multilineaire

In [47]:

```
#Multilinear
#extract x and y from our data
x=house[["sqft_living","grade"]]
  #we have more than one input
y=house["price"].values
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.35,random_state=40) #splitting data
with test size of 35%


model=LinearRegression() #build linear regression model
model.fit(x_train,y_train) #fitting the training data
predicted=model.predict(x_test) #testing our model's performance

print("MSE", mean_squared_error(y_test,predicted))
print("R squared", metrics.r2_score(y_test,predicted))
```
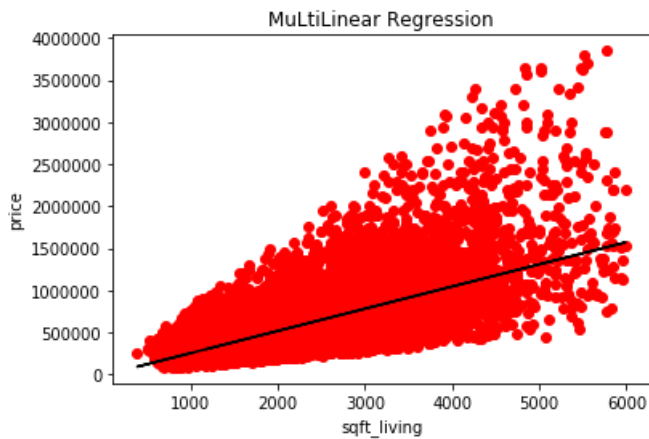
```
MSE 53467526289.49688
R squared 0.527940052513393
```

In [45]:

```
plt.scatter(x, y, color="r")
plt.title("MuLtiLinear Regression")
plt.ylabel("price")
plt.xlabel("sqft_living")
plt.plot(x,model.predict(x),color="k")
plt.show()
```

```python
plt.show()
```



Regression Polynomiale

```python
#Importing libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import PolynomialFeatures

from sklearn.metrics import mean_squared_error
from sklearn import metrics
x= house[["sqft_living", "grade"]]
y= house["price"].values

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.35, random_state=40)  #splitt
ing data
lg=LinearRegression()
poly=PolynomialFeatures(degree=3)

x_train_fit = poly.fit_transform(x_train) #transforming our input data
lg.fit(x_train_fit, y_train)
x_test_ = poly.fit_transform(x_test)
predicted = lg.predict(x_test_)

print("MSE: ", metrics.mean_squared_error(y_test, predicted))
print("R squared: ", metrics.r2_score(y_test,predicted))
```

```
MSE:  49231222144.91536
R squared:  0.5647471380908822
```
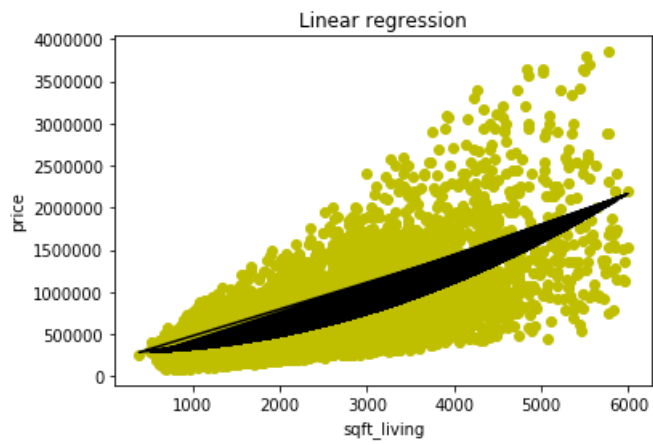
In [79]:

```python
x= house[["sqft_living"]]
y= house["price"].values.reshape(-1,1)
poly = PolynomialFeatures(degree = 2)
x_poly = poly.fit_transform(x)
poly.fit(x_poly, y)
lg=LinearRegression()
lg.fit(x_poly, y)

plt.scatter(x, y, color="y")
plt.title("Linear regression")
plt.ylabel("price")
plt.xlabel("sqft_living")
plt.plot(x, lg.predict(poly.fit_transform(x)), color="k")
```

Out[79]:

```
[<matplotlib.lines.Line2D at 0x12ca3208>]
```

Linear regression

In [59]:

```
print(x.shape)
```

(21268, 1)

In [60]:

```
print(y.shape)
```

(21268, 1)

In [ ]: